

Databases II: DB programming

Datawarehousing
&
Business Intelligence

What do we learn?

- basic concepts, advantages and disadvantages of Data Warehousing
- differences between OLTP and Data Warehousing systems
- architecture and most important components of a Data Warehouse
- how Data Warehousing has evolved
- problems with Data Warehousing
- the concept of a data mart
- activities to start a Data Warehouse project
- Kimball's methodology
- concepts of Dimensional modelling
- creation of a DWH

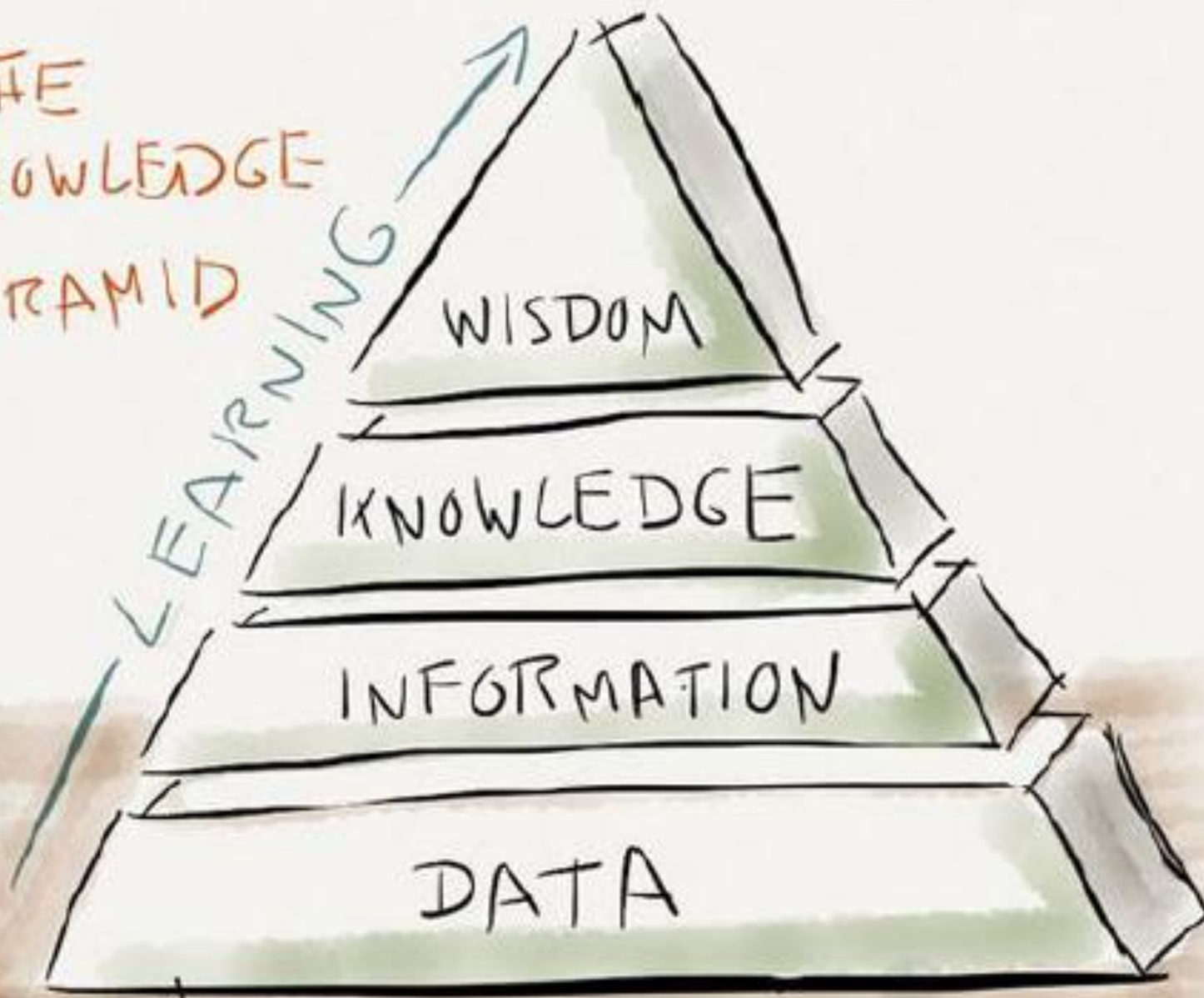
Introduction

Introduction

There is a growing need to turn data into information

- flexible business reporting available for the business
- data in DWH grows exponentially
 - terabytes of data in a DWH are 'normal'
- applications using data have become more complex
 - traditional reporting
 - advanced analysis
- all traditional DBMS offer DWH facilities

THE KNOWLEDGE PYRAMID



Business Intelligence (BI): definition

Business Intelligence(BI) comprises the set of **strategies, processes, applications, data, technologies** and **technical architectures** which are used by enterprises to support the collection, data analysis, presentation and dissemination of business information.

BI technologies provide **historical, current and predictive views of business operations**.

Common functions of business intelligence technologies include **reporting, online analytical processing, analytics, data mining, process mining, complex event processing, business performance management, benchmarking, text mining, predictive analytics** and **prescriptive analytics**.

Source: wikipedia

Business Intelligence (BI): another definition

Business intelligence, or BI, is a blanket term for technology that helps companies organize and query their data to help them improve operations and make more money.

Source: <http://fortune.com/2017/04/25/infor-buys-business-intelligence-player>

Drivers for increasing use of BI

- Digitalisation (ERP, CRM, PLM, DAM, PIM, ...)
- Connectors between BI software and Business software
 - Way of working can stay the same
 - BI comes on top of the existing SW
- Data overflow
- Complexity and Speed of change in the Business Environment
 - Gut feeling and experience is often not sufficient anymore
- Reduce inefficiencies, inaccuracies
 - anti –Excel culture
- Decreasing Cost

BI technology vendors

- Microsoft:
 - Reporting: Microsoft reporting services + PowerBI
 - BI+data mining: SSAS (SQL Server Analysis Services)
 - ETL: SSIS (SQL Server Integration Services)
- Cognos (now IBM)
 - ETL
 - Reporting tools
- Business Objects (now SAP)
 - Reporting
 - ETL
- SAP Business Warehouse
 - Kubus
- Tableau
- Datastage ETL
- QlikView reporting

Gartner Magic Quadrant for Business Intelligence and Analytics Platforms, February 2018

Figure 1. Magic Quadrant for Analytics and Business Intelligence Platforms



Source: Gartner (February 2018)

Datawarehouse: definition

a **data warehouse** is an integrated, subject oriented, time variant and non volatile collection of data to support decisions taken on management level.

- subject oriented
 - The warehouse is organized around the major subjects of the enterprise (e.g. customers, products, and sales) rather than the major application areas (e.g. customer invoicing, stock control, and product sales).
 - This is reflected in the need to store decision-support data rather than application-oriented data.

Properties

- integrated
 - The data warehouse integrates corporate application-oriented data from different source systems, which often includes data that is inconsistent.
- time
 - Data in the warehouse is only accurate and valid at some point in time or over some time interval.
 - Time-variance is also shown in the extended time that the data is held, the implicit or explicit association of time with all data, and the fact that the data represents a series of snapshots.
 - e.g. unit price of a product can be stored historically or can be a snapshot.
 - DW will build history by taking regular snapshots.

Properties

- non volatile
 - Data in the warehouse is not normally updated in real-time (RT) but is refreshed from operational systems on a regular basis. (However, emerging trend is towards RT or near RT DWs)
 - New data is always added as a supplement to the database, rather than a replacement.
- aggregated data
 - e.g. generated by GROUP BY

Goals of DWH

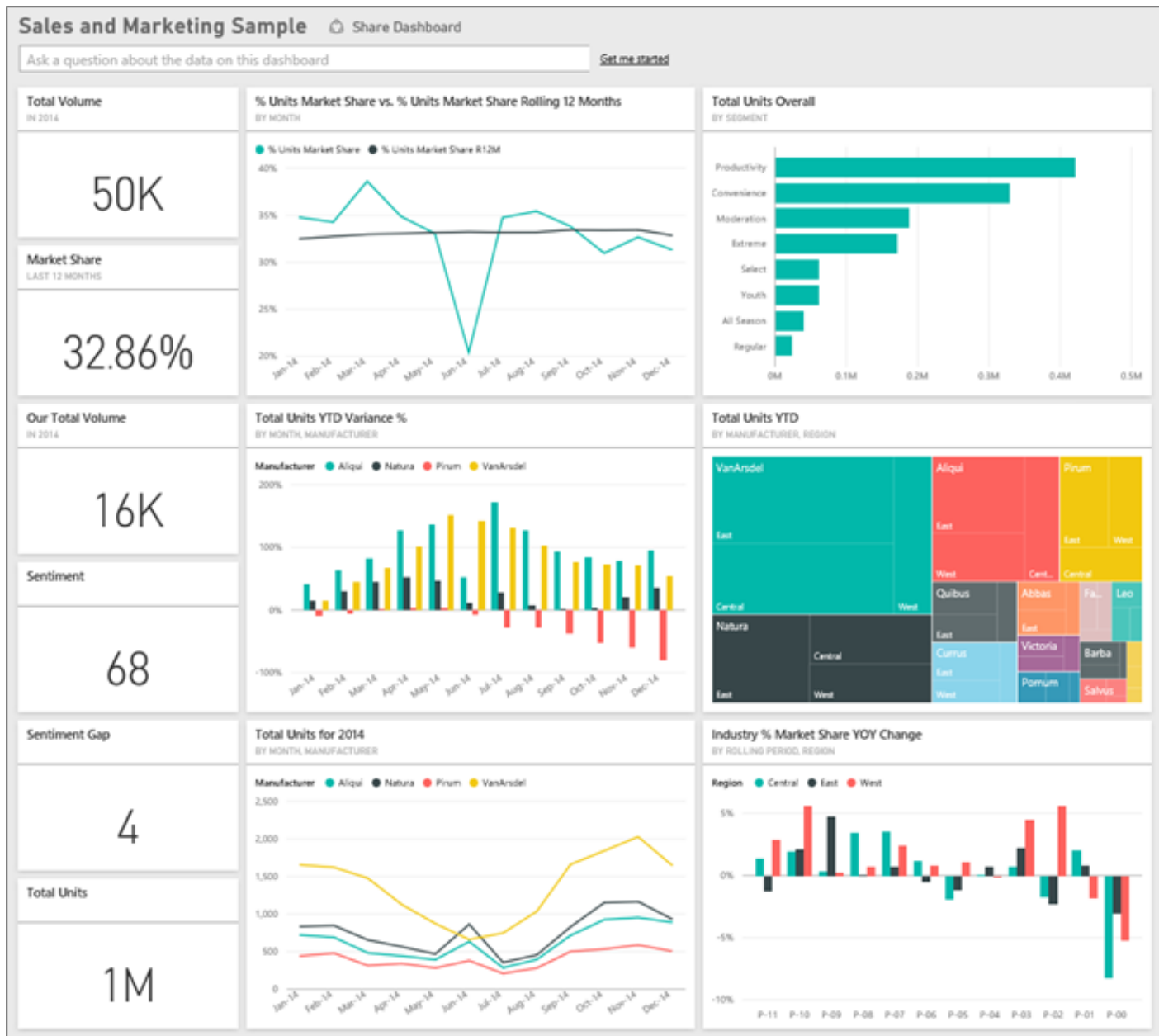
- Reporting
- Analysis of events in the past or actual events
- Predictions based on trend analysis
- Multidimensional reporting
- Empowerment of end user by offering simplified reporting tools (cf. SQL: only specialist can write SQL)
- Data mining

DWH is a technology used to realize BI solutions,
but it is by far not the only technology.

e.g. multidimensional report

Genre	2013	2014	2015	2016	2017	Total
Alternative		5.94	3.96	3.96		13.86
Alternative & Punk	62.37	39.60	45.54	38.61	55.44	241.56
Blues	10.89	10.89	19.80	8.91	9.90	60.39
Bossa Nova	0.99	1.98	7.92		3.96	14.85
Classical		13.86	9.90	16.83		40.59
Comedy		3.98	1.99	11.94		17.91
Drama		17.91	11.94	17.91	9.95	57.71
Easy Listening	2.97	1.98	2.97		1.98	9.90
Electronica/Dance	1.98	3.96	1.98	2.97	0.99	11.88
Heavy Metal	3.96	2.97		2.97	1.98	11.88
Hip Hop/Rap	1.98	2.97	3.96	3.96	3.96	16.83
Jazz	19.80	15.84	15.84	5.94	21.78	79.20
Latin	82.17	77.22	80.19	63.36	79.20	382.14
Metal	61.38	53.46	25.74	65.34	55.44	261.36
Pop	0.99	8.91	9.90	7.92		27.72
R&B/Soul	7.92	8.91	6.93	9.90	6.93	40.59
Reggae	5.94	6.93	7.92	5.94	2.97	29.70
Rock	178.20	155.43	156.42	162.36	174.24	826.65
Rock And Roll	0.99	1.98	0.99	1.98		5.94
Sci Fi & Fantasy		9.95	17.91	11.94		39.80
Science Fiction		3.98	3.98	1.99	1.99	11.94
Soundtrack	3.96	3.96	4.95	3.96	2.97	19.80
TV Shows		25.87	27.86	25.87	13.93	93.53
World	2.97	2.97	0.99	2.97	2.97	12.87
Total	449.46	481.45	469.58	477.53	450.58	2,328.60

e.g. BI dashboard for manager



Advantages

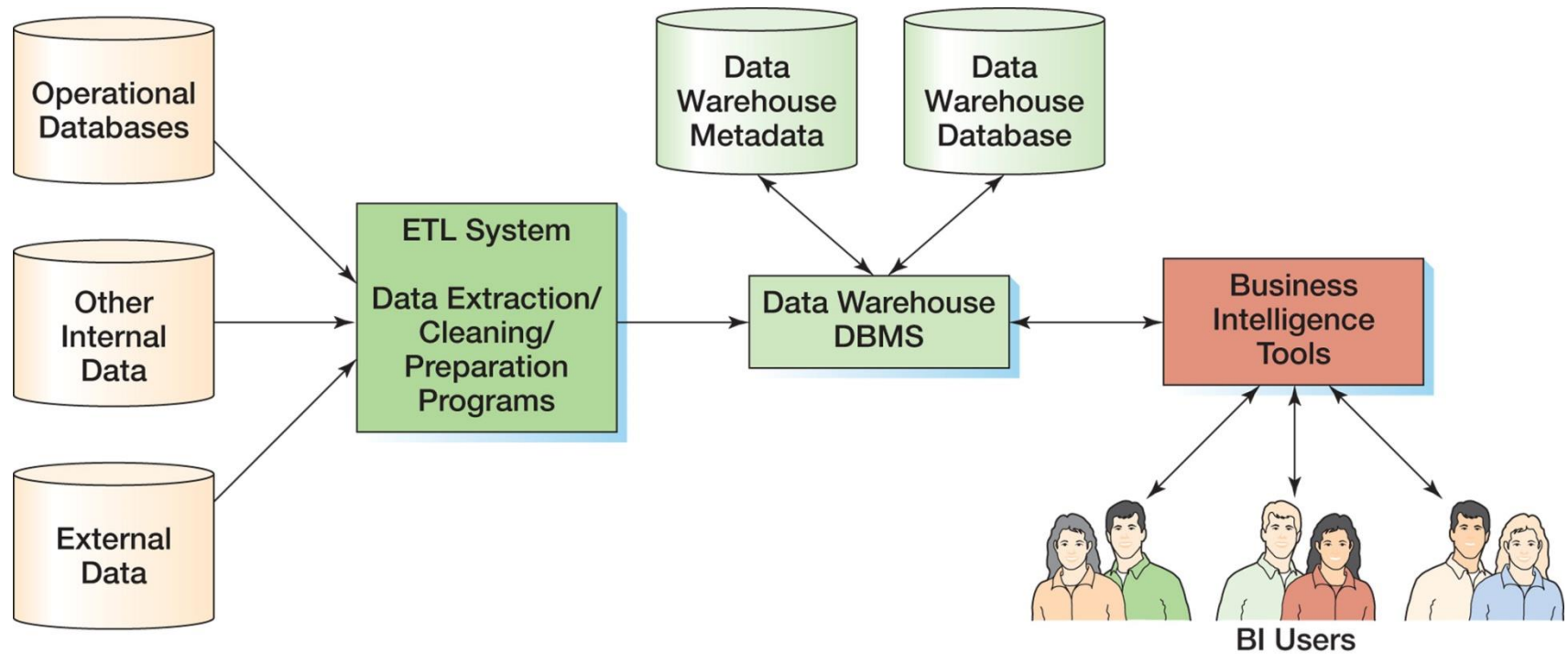
- high ROI (return on investment)
 - large investments with potentials high ROI after a short period
- Competitive advantage
 - decision makers get access to data that was not available, unknown or unused before.
- Increased productivity of corporate decision-makers
 - decision maker gets one consistent view on the enterprise
 - because data of different sources can be integrated to one consistent view, that is subject oriented and keeps history.
 - decision makers can make more substantial, more accurate and more consistent analysis.
 - tools can help turn data into useful information

Comparison of OLTP Systems and Data Warehousing

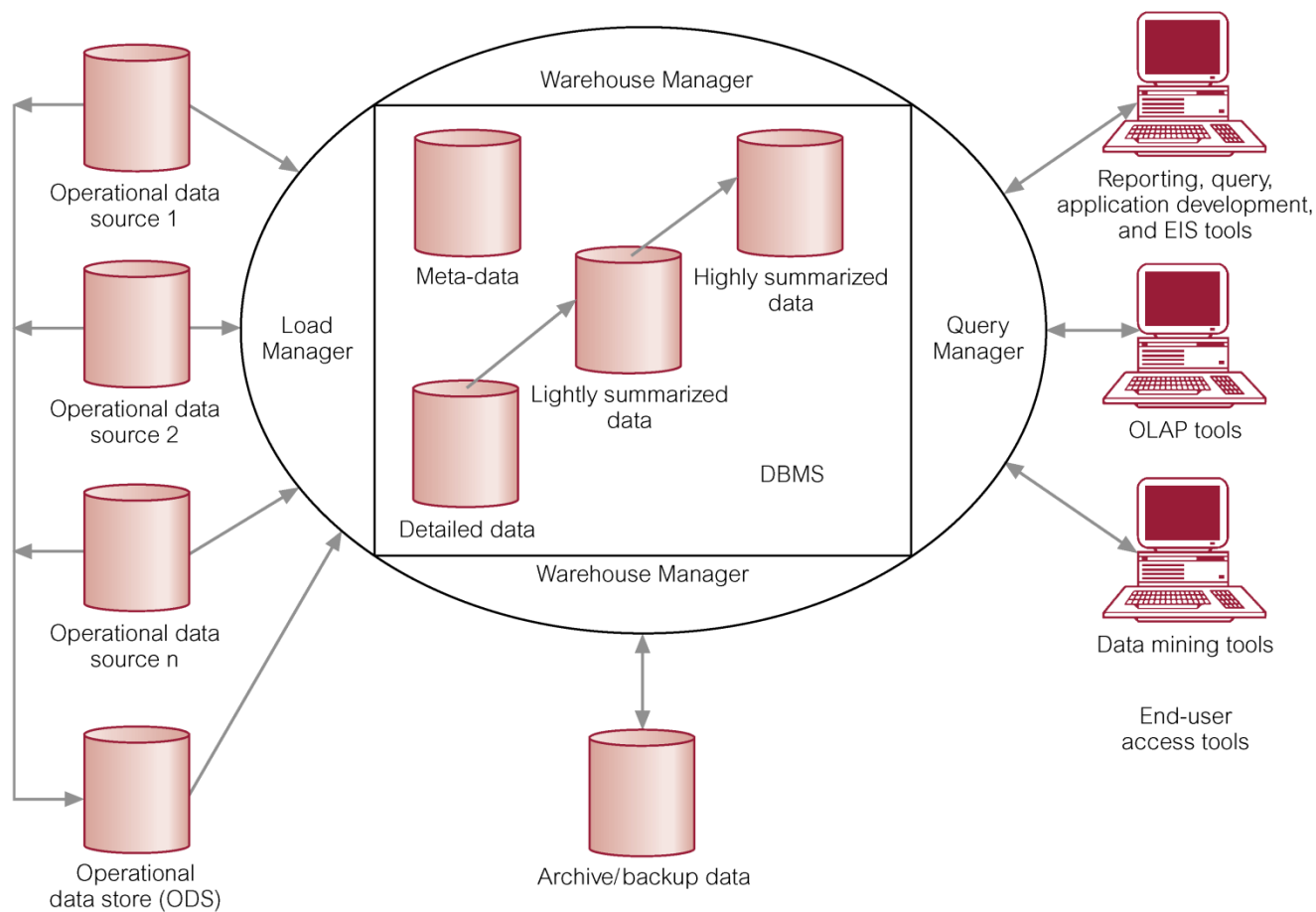
CHARACTERISTIC	OLTP SYSTEMS	DATA WAREHOUSING SYSTEMS
Main purpose	Support operational processing	Support analytical processing
Data age	Current	Historic (but trend is toward also including current data)
Data latency	Real-time	Depends on length of cycle for data supplements to warehouse (but trend is toward real-time supplements)
Data granularity	Detailed data	Detailed data, lightly and highly summarized data
Data processing	Predictable pattern of data insertions, deletions, updates, and queries. High level of transaction throughput.	Less predictable pattern of data queries; medium to low level of transaction throughput
Reporting	Predictable, one-dimensional, relatively static fixed reporting	Unpredictable, multidimensional, dynamic reporting
Users	Serves large number of operational users	Serves lower number of managerial users (but trend is also toward supporting analytical requirements of operational users)

Architecture

DWH components

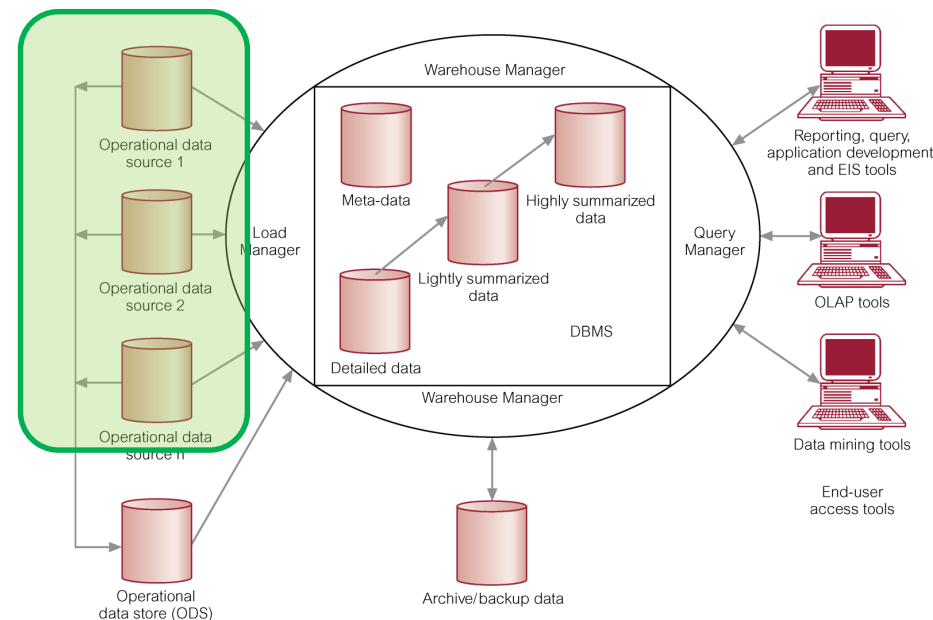


Architecture of a DWH



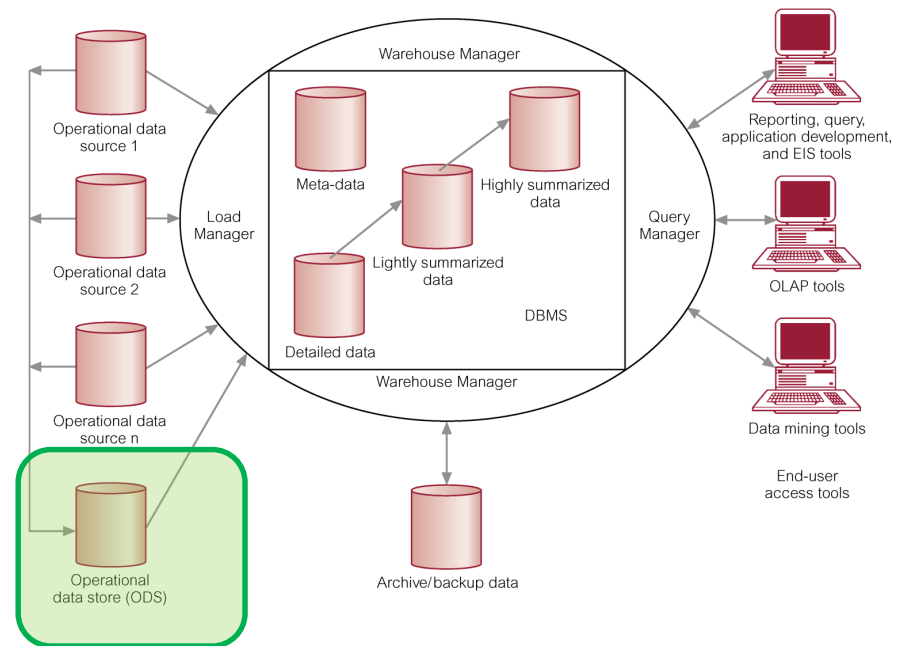
Architecture of a DWH

- operational data
- sources of data
 - mainframe (hierarchical, network, relational db)
 - departmental data in files and RDBM systems
 - private data on workstations and private servers
 - external systems
 - internet, e.g. prices of competitors
 - commercial DB
 - data used by customers or suppliers



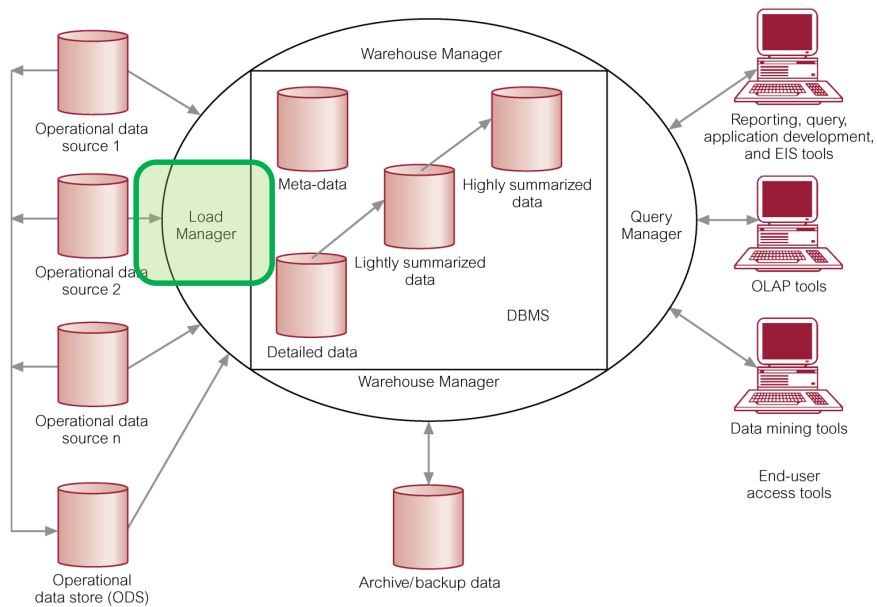
Architecture of a DWH

- operational data source
- repository
 - current, integrated data
 - preparing step in development of DWH,
- or
 - support for reporting services in legacy systems



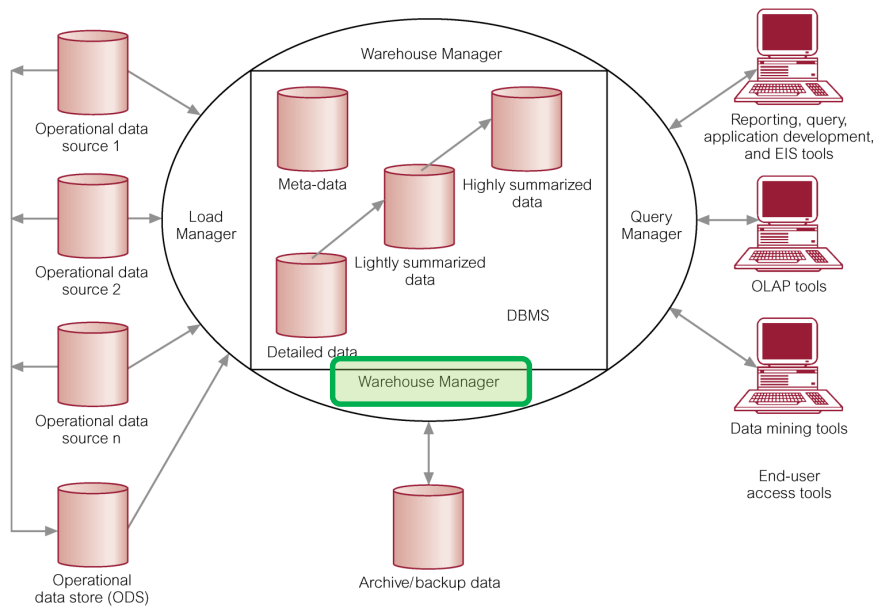
Architecture of a DWH

- ETL manager
 - supports all ETL operations (Extraction, Transformation and Load) on data
 - either directly on operational data
 - or on operational data store



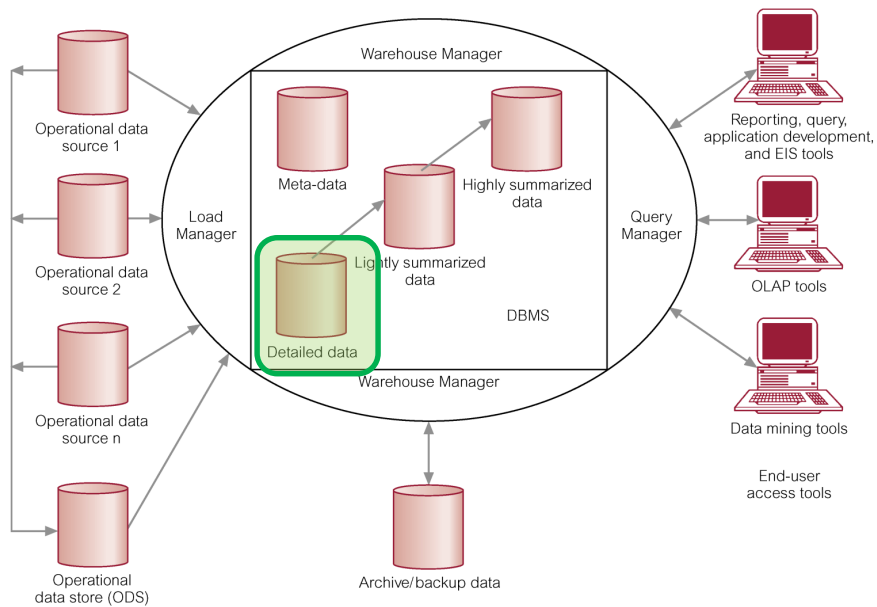
Architecture of a DWH

- warehouse manager
 - management of data in the DWH
 - analysis of data to guarantee consistency
 - transformation and merging data from sources or temporary storage in DWH tables
 - creation of indexes and views
 - possibly denormalised
 - creation of aggregates (data summary)
 - back-up and archiving of the data



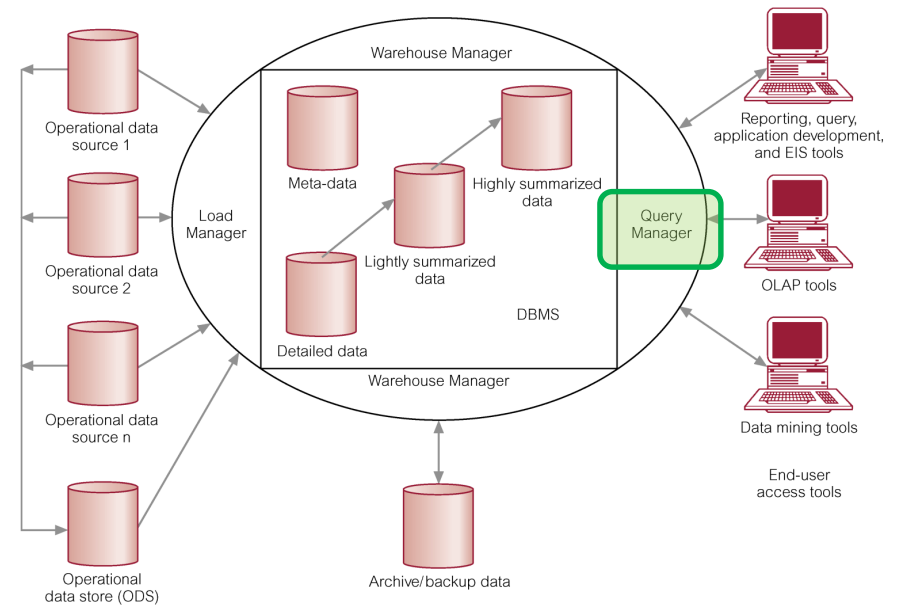
Architecture of a DWH

- detailed data
 - this data is added to the DWH on a regular (e.g. daily) basis



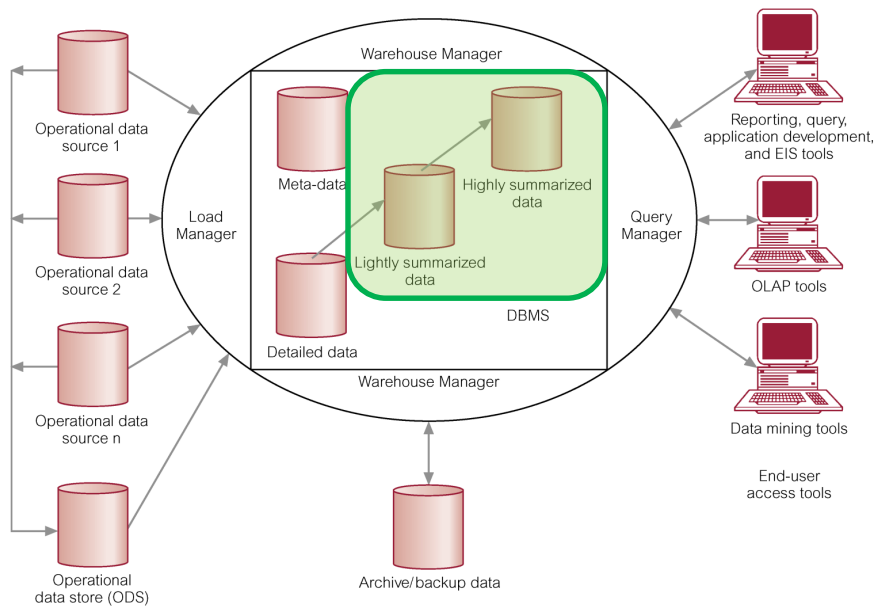
Architecture of a DWH

- query manager
 - management of user queries
 - use of correct tables
 - execution/scheduling of queries
 - profile generations
 - proposals for aggregates and indexes



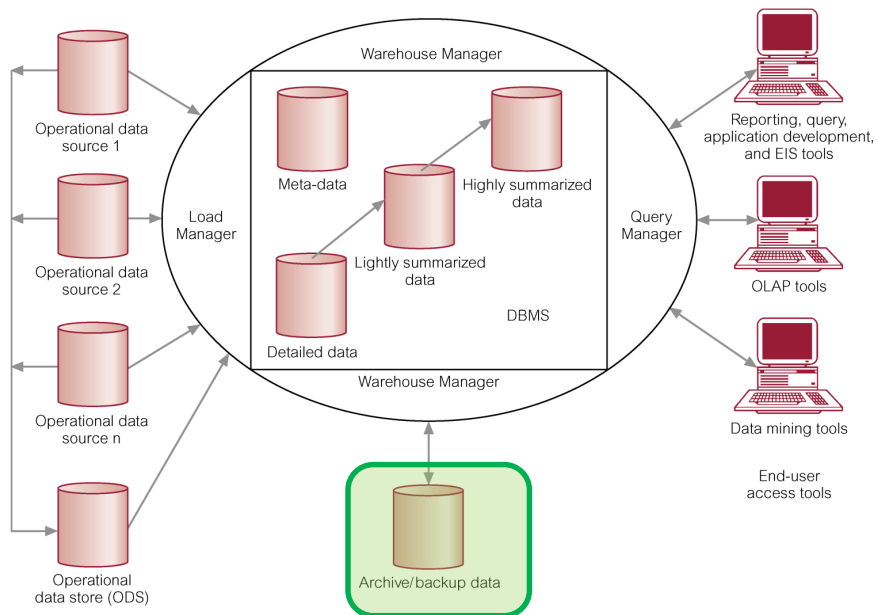
Architecture of a DWH

- summarised data
 - predefined summarised data
 - can be adapted in a flexible way so different sorts of queries are supported
 - improved query execution performance as opposed to detailed data because data is "prepared"



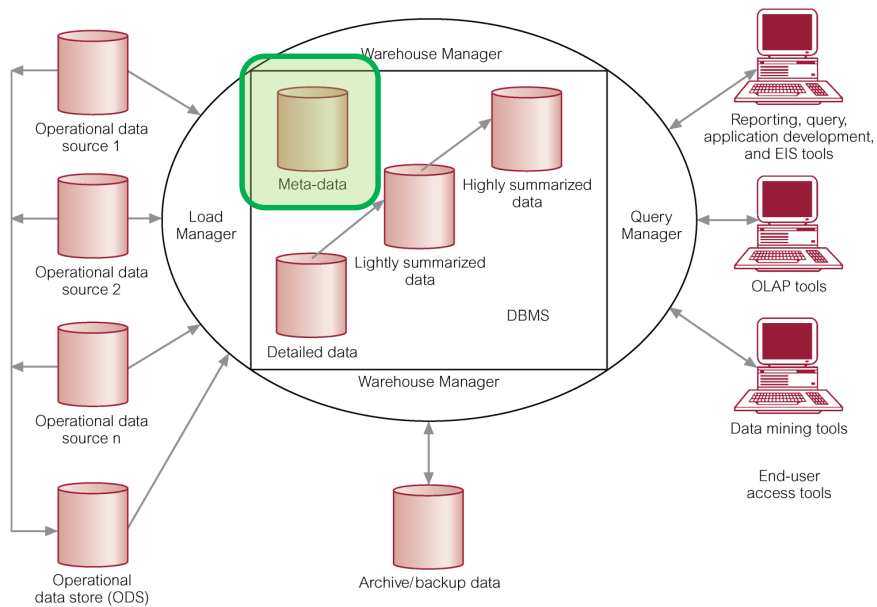
Architecture of a DWH

- archive/back-up data
 - for both detailed and backup data
 - summarised data can be kept longer than detailed data



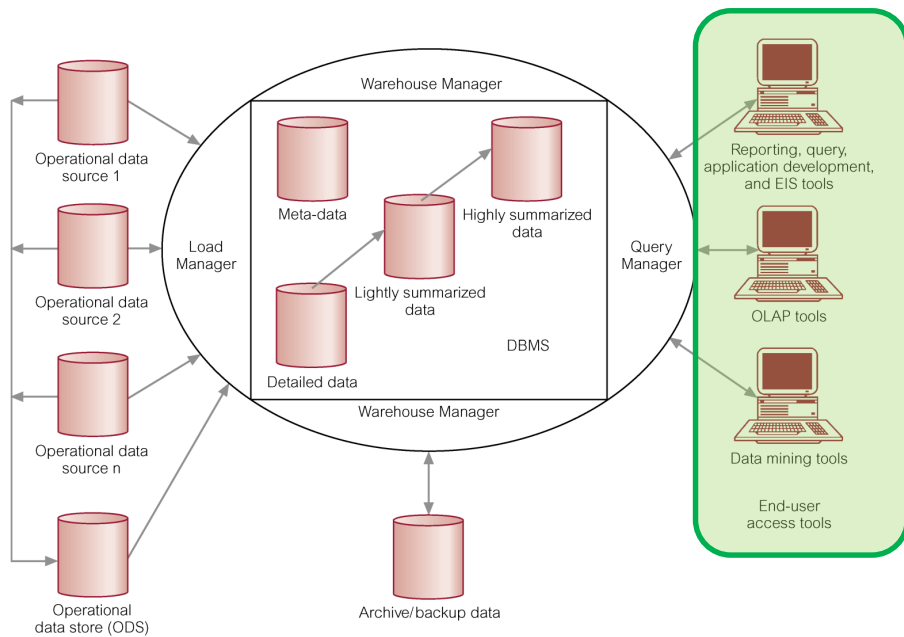
Architecture of a DWH

- meta data
 - necessary for
 - ETL
 - DWH manager
 - Query manager
- several versions of metadata each adapted to a specific process
- allows to determine the source of each data item in the DWH

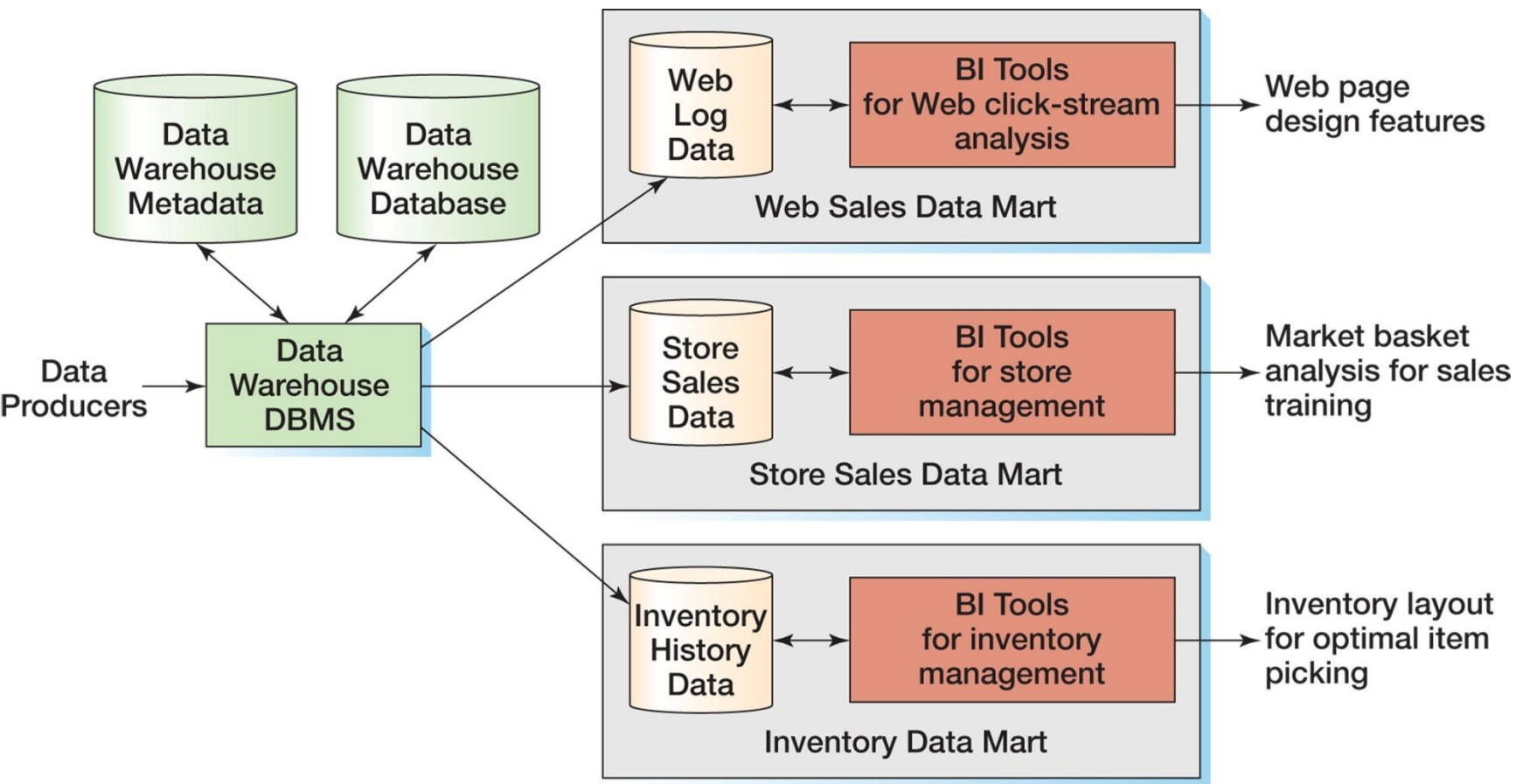


Architecture of a DWH

- end user access tools
 - reporting and querying
 - application development tools
 - OLAP tools
 - data mining tools



DWH & Data marts



Datamart

a DB existing of a **subset of company data** to support the needs of a particular business unit for data analysis, or, to support users sharing the same needs for analyzing **business processes**.

- Why a data mart?
 - to give users access to the data they analyse most frequently
 - to offer data in a way that corresponds to the collective view of a group of users in a department or a group of users in the same business process
 - to improve response time by offering lower data volumes
 - to offer data in a format that fits the tools used by end users (OLAP, datamining tools)
 - reduction of complexity in the ETL process
 - reduction of cost as opposed to the setup of a complete DWH

Data Mining Applications

- **Data mining** are use to:
 - perform what-if analyses
 - perform predictions ("predictive analysis")
 - facilitate the decision process
- Data mining applications use sophisticated statistical and mathematical techniques
- Reports are less critical

Problems associated with DWH

- **Underestimation of resources (costs) for ETL**
 - extraction, transformation and loading of data in the DWH takes a major part of the development time
 - Projects might take years
- **Hidden problems with source systems**
 - are sometimes only discovered after years
 - can be solved either in DWH or in operational DB
 - e.g. fields allowing NULL values: some offices never fill out the fields, although data is available and can be useful.
- **Required data not captured**
 - change actual system or develop separate system for those data
 - e.g. the data a customer has registered is not captured
- **Increased end-user demands**
 - enhanced load of IT personnel
 - demand for more user friendly, powerful and sophisticated tools
 - better end user training

Problems associated with DWH

- **Data homogenization**
 - trying to focus on similarities between data might lower the use of the data
 - e.g. similarities between sale and rent of properties
- **Need for concurrent support of several (historical) versions**
 - Operational systems evolve (i.e. database schema changes over time) but data from older versions and newer versions resides together in the DWH
 - Might be challenging
- **High demand for resources**
 - e.g. disk space
- **Data ownership**
 - sensitive departmental data (e.g. HR) is highly secured in HR department but is widely available in DWH
- **High maintenance**
 - each change in business processes or in sources systems influences the DWH (both structure and ETL)

Problems associated with DWH

- **Long projects**
 - development can take years
 - phased development through data marts (see further) might be recommended
- **DWH creates expectation of user ‘empowering’:**
 - make own reports, analyses
 - less need for IT
 - but:
 - meta dictionary that describes data in DWH is necessary
 - dependency of a few specialists remains
- **complexity of integration**
 - different DWH tools have to work together smoothly
- **Complex change and version management**
 - Consistency in reporting between versions of underlying databases
- **Night might be too short for ETL**

Problems with Operational data

- Dirty data
- Missing values
- Inconsistent data
- Data not integrated
- Wrong format
 - Too fine
 - Not fine enough
- Too much data
 - Too many attributes
 - Too much volume

Design

Design of a DWH

- 2 development methodologies
- **Inmon**
 - creation of a data model based on all data of the organisation
 - Enterprise Data Warehouse (EDW)
 - Used to distil data marts for each department
 - uses traditional methods for describing EDW:
 - ERD
 - tables in normal form
- **Kimball**
 - Starts by identifying the information requirements (referred to as analytical themes) and associated business processes of the enterprise.
→ Data Warehouse Bus. Matrix
 - This first data mart is critical in setting the scene for the later integration of other data marts as they come online.
 - The integration of data marts ultimately leads to the development of an EDW.
 - Uses dimensionality modelling to establish the data model (called star schema) for each data mart.

Kimball's Business Dimensional Lifecycle

- Guiding principle
 - meet the information requirements of the enterprise by building
 - single,
 - integrated,
 - easy-to-use,
 - high-performanceinformation infrastructure, which is delivered in meaningful increments of 6 to 12 month timeframes.
- Goal
 - deliver the entire solution including
 - the data warehouse,
 - *ad hoc* query tools,
 - reporting applications,
 - advanced analytics
 - all the necessary training and support for the users.

Data Warehouse Sample database

- **Adventureworks Database**
 - Sample Database used by many books and samples about MS SQL Server
 - Exists in an OLTP and DW version
 - Mainly stores data about Products and Sales
 - Download backup files from <https://github.com/Microsoft/sql-server-samples/releases/tag/adventureworks>

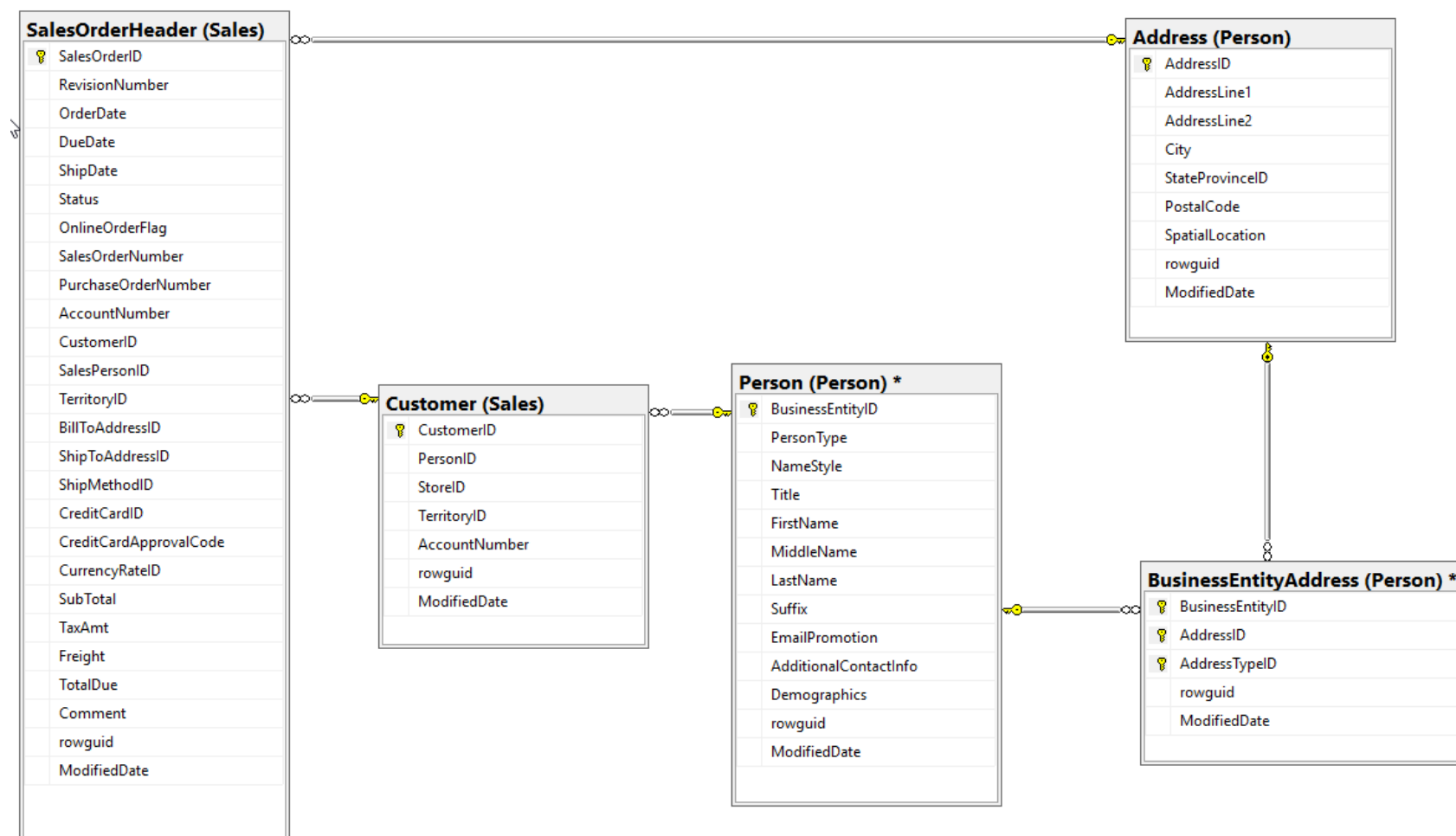
Data Warehouse Sample Queries

- What was the total revenue for UK in the third quarter of 2010?
- What was the total revenue for bike sales for each type of bike (mountain or road) in Germany in 2011?
- What are the three most popular bike types 2013 and how does this compare with the figures for the previous two years?
- What is the monthly revenue for clothing in each region, compared with rolling 12-monthly prior figures?
- What is the relationship between the total annual revenue for bikes and the average temperature for each country?

Dimensionality modelling: example

- ERD AdventureWorks

Often data in an OLTP system is highly normalized, which causes you might need all these tables for a simple sales report:



Dimensionality modelling: example

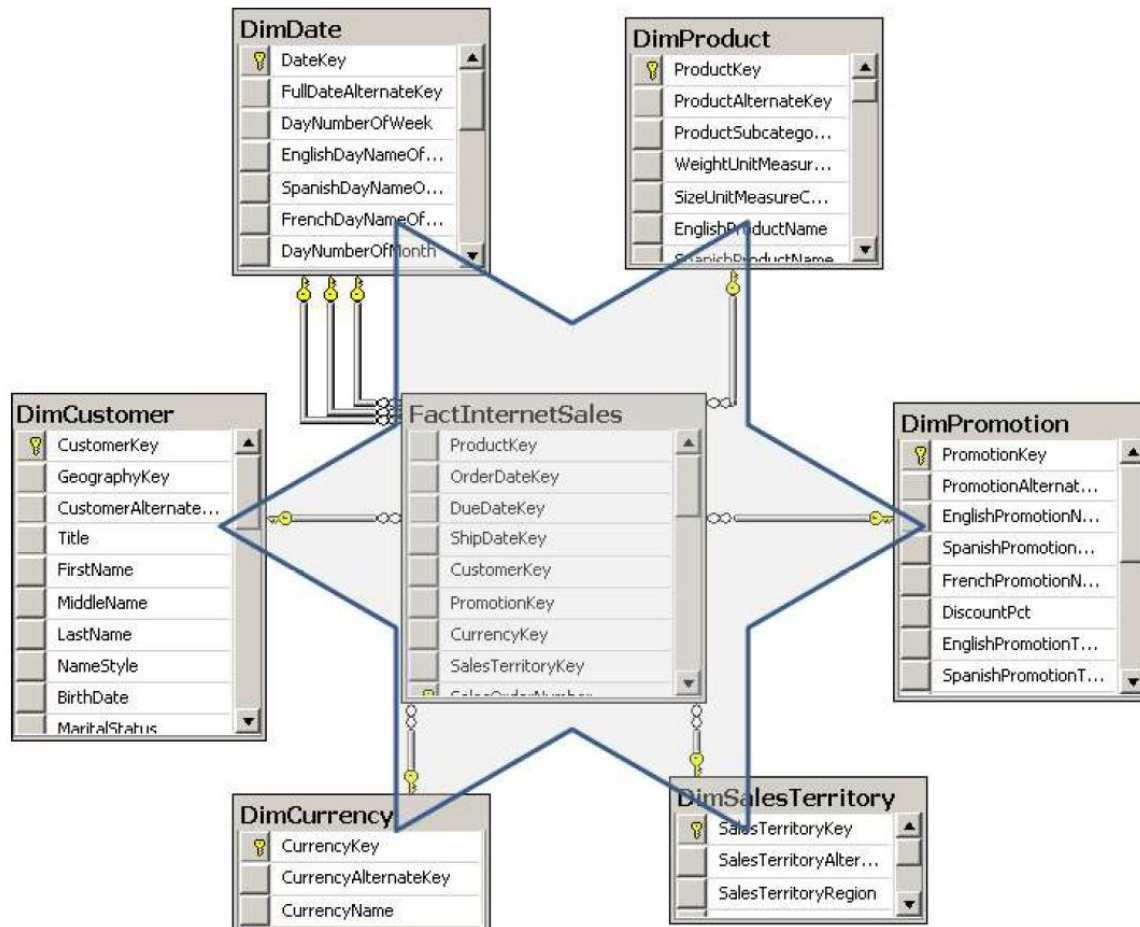
• ERD AdventureWorks

The query for a simple sales report might look like this (total sales amount per year, quarter and city):

```
select year(orderdate) as calenderyear,
(month(orderdate)-1)/3+1 as
calenderquarter, pa.City,
sum(subtotal)
from sales.salesorderheader s
join sales.customer c on
s.CustomerID=c.CustomerID
join person.person p on
c.PersonID=p.BusinessEntityID
join person.BusinessEntityAddress a on
p.BusinessEntityID=a.BusinessEntityID
join person.address pa on
a.AddressID=pa.AddressID
group by year(orderdate), (month(orderdate)-
1)/3+1, pa.City
order by 1,2,3;
```

calenderyear	calenderquarter	City	(No column name)
2011	2	Ballard	3578,27
2011	2	Bellflower	699,0982
2011	2	Bellingham	3578,27
2011	2	Bendigo	20909,78
2011	2	Berkeley	3578,27
2011	2	Berlin	3399,99
2011	2	Beverly Hills	699,0982
2011	2	Birmingham	3399,99
2011	2	Bremerton	4277,3682
2011	2	Brisbane	3578,27
2011	2	Burbank	7156,54
2011	2	Caloundra	3578,27
2011	2	Coffs Harbour	699,0982
2011	2	Colombes	3374,99
2011	2	Concord	3578,27
2011	2	Coronado	7156,54
2011	2	Courbevoie	3374,99
2011	2	Cranbourne	3578,27
2011	2	Daly City	3578,27
2011	2	Dedham	3578,27

Star schema



A Star schema is a logical structure that has **a fact table** (containing factual data) in the center, **surrounded by denormalized dimension tables** (containing reference data).

• remark:

- natural keys from the operational system are available but not as a key in the star schema
- surrogate integer keys are used instead because they are simpler and faster
- this way independence between OLTP and DWH is ensured

Star schema

To get the same results on the star schema the query looks like this:

```
SELECT calendaryear, calendarquarter, City, sum(salesamount)
FROM factinternetsales AS f
    JOIN
    dimdate AS d
    ON f.orderdatekey = d.datekey
    JOIN
    dimcustomer AS c
    ON f.CustomerKey = c.CustomerKey
    JOIN
    dimgeography AS g
    ON c.GeographyKey = g.GeographyKey
GROUP BY calendaryear, calendarquarter, city
ORDER BY calendaryear, calendarquarter, city;
```

- conclusion: the query on the OLTP system
 - has deeper join constructs
 - joins more tables
 - performs calculations (e.g. to get the quarter)



star schema

- **fact table** contains data about facts
 - e.g. factual data about sales of property: sales price, commission percentage, ...
- **dimension table** contains reference information
 - e.g. property data (address, etc), buyer, owner, ...
- facts are generated by events that happened (e.g. a sale)
- most probably facts never change

the fact table





- Bulk of data in data warehouse is in fact tables, which can be extremely large.
- Important to treat fact data as read-only reference data that will not change over time.
- Most useful fact tables contain one or more numerical measures, or '**facts**' that occur for each record and are numeric and additive.


FactInternetSales

ProductKey
OrderDateKey
DueDateKey
ShipDateKey
CustomerKey
PromotionKey
CurrencyKey
SalesTerritoryKey
 SalesOrderNumber
 SalesOrderLineNumber
RevisionNumber
OrderQuantity
UnitPrice
ExtendedAmount
UnitPriceDiscountPct
DiscountAmount
ProductStandardCost
TotalProductCost
SalesAmount
TaxAmt
Freight
CarrierTrackingNumber
CustomerPONumber
OrderDate
DueDate
ShipDate

the dimension tables

- Dimension tables usually contain descriptive textual information.
- Dimension attributes are used as the constraints in data warehouse queries. (e.g. in where or having clauses)
- Star schemas can be used to speed up query performance by denormalizing reference information into a single dimension table.

DimDate	
 DateKey	
FullDateAlternateKey	
DayNumberOfWeek	
EnglishDayNameOfWeek	
SpanishDayNameOfWeek	
FrenchDayNameOfWeek	
DayNumberOfMonth	
DayNumberOfYear	
WeekNumberOfYear	
EnglishMonthName	
SpanishMonthName	
FrenchMonthName	
MonthNumberOfYear	
CalendarQuarter	
CalendarYear	
CalendarSemester	
FiscalQuarter	
FiscalYear	
FiscalSemester	

DimSalesTerritory	
 SalesTerritoryKey	
SalesTerritoryAlternateKey	
SalesTerritoryRegion	
SalesTerritoryCountry	
SalesTerritoryGroup	
SalesTerritoryImage	

Snowflake schema

DimDate	
⚡	DateKey
	FullDateAlternateKey
	DayNumberOfWeek
	EnglishDayNameOfWeek
	SpanishDayNameOfWeek
	FrenchDayNameOfWeek
	DayNumberOfMonth
	DayNumberOfYear
	WeekNumberOfYear
	EnglishMonthName
	SpanishMonthName
	FrenchMonthName
	MonthNumberOfYear
	CalendarQuarter
	CalendarYear
	CalendarSemester
	FiscalQuarter
	FiscalYear
	FiscalSemester

FIGURE 1-4 The *DimDate* denormalized dimension.

Snowflake schema is a variant of the star schema that has a **fact table in the centre**, surrounded by **normalised dimension tables**.

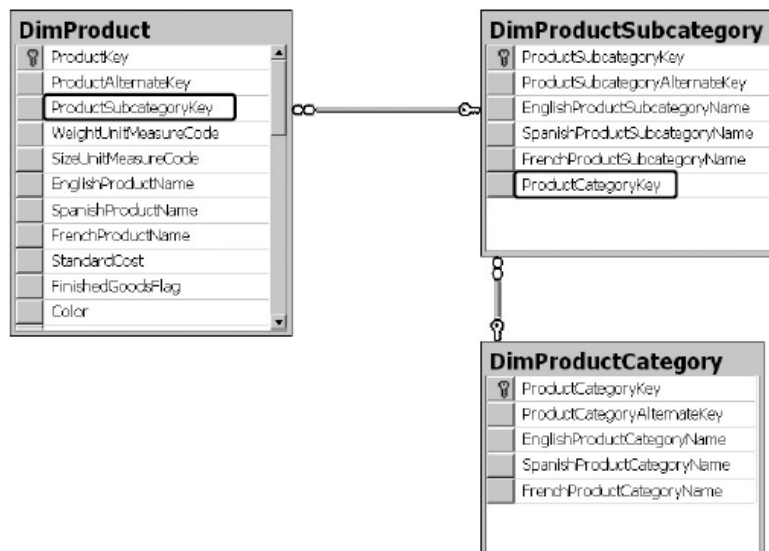



FIGURE 1-5 The *DimProduct* normalized dimension.

Slowly Changing Dimensions

- = Dimensions that change slowly and irregularly over a period of time
- We want to keep track of old and current values of these dimensions
- Example: product information can change over time, e.g. the color of a ProductID can change.
- The most common way to handle this is adding two date fields (Start and End) to hold the validity period. End = NULL means currently valid

	Column Name	Data Type	Allow Nulls
	ProductKey	int	<input type="checkbox"/>
	ProductID	int	<input type="checkbox"/>
	Name	nvarchar(50)	<input checked="" type="checkbox"/>
	Color	nvarchar(15)	<input checked="" type="checkbox"/>
	ListPrice	money	<input checked="" type="checkbox"/>
	Size	nvarchar(50)	<input checked="" type="checkbox"/>
	Weight	decimal(8, 2)	<input checked="" type="checkbox"/>
	Start	date	<input checked="" type="checkbox"/>
	[End]	date	<input checked="" type="checkbox"/>

```
select * from dimproduct where productid = 776;
```

161 %

Results Messages

	ProductKey	ProductID	Name	Color	ListPrice	Size	Weight	Start	End
1	281	776	Mountain-100 Black, 42	Black	3374,99	42	20.77	2011-05-31	2019-01-23
2	506	776	Mountain-100 Black, 42	Blue	3374,99	42	20.77	2019-01-23	NULL

advantages of the dimensional model

- Predictable and standard form of the underlying dimensional model offers important advantages:
 - **Efficiency**
 - a consistent DB structure allows tools to have efficient access to the data
 - **Ability to handle changing requirements**
 - the model can easily adapt to changing needs because each dimension is equivalent to the fact table
 - ideal for ad hoc queries
 - **Extensibility**
 - adding new facts
 - adding new dimensions
 - adding attributes to dimensions
 - **Ability to model common business situations**
 - **Predictable query processing**
 - the way tables are used is predictable (not the queries themselves)

DM and ER models

- Entity Relationship Diagrams
 - used to design the DB of OLTP systems
 - basics: remove redundancies
 - redundancy causes update/delete/insert anomalies
 - ad hoc queries are more difficult
 - lots of tables can be involved: deep join constructs
- Dimensional Modelling
 - used for design of DWH of data mart
 - intuitive storage and high performance consulting of data
- A single ER model normally decomposes into multiple DMs.
- Multiple DMs are then associated through ‘shared’ dimension tables.

Dimensional modelling Stage

- Design issues
 - Choose granularity level:
 - type 1: the number of dimensions determines the granularity level of the analysis you can get
 - type 2: every order, summarized by month, quarter, ...
 - Duration measures how far back in time the fact table goes.
 - Slowly changing dimension problem means that the proper description of the old dimension data must be used with the old fact data.
 - type 1: the attribute that changes is **overwritten**
 - type 2: if an attribute changes a **new record** is added to the dimension table
 - type 3: make sure **old and new value** of the attribute is available in the same record

Exercise

- See document
“SSIS sample Project”